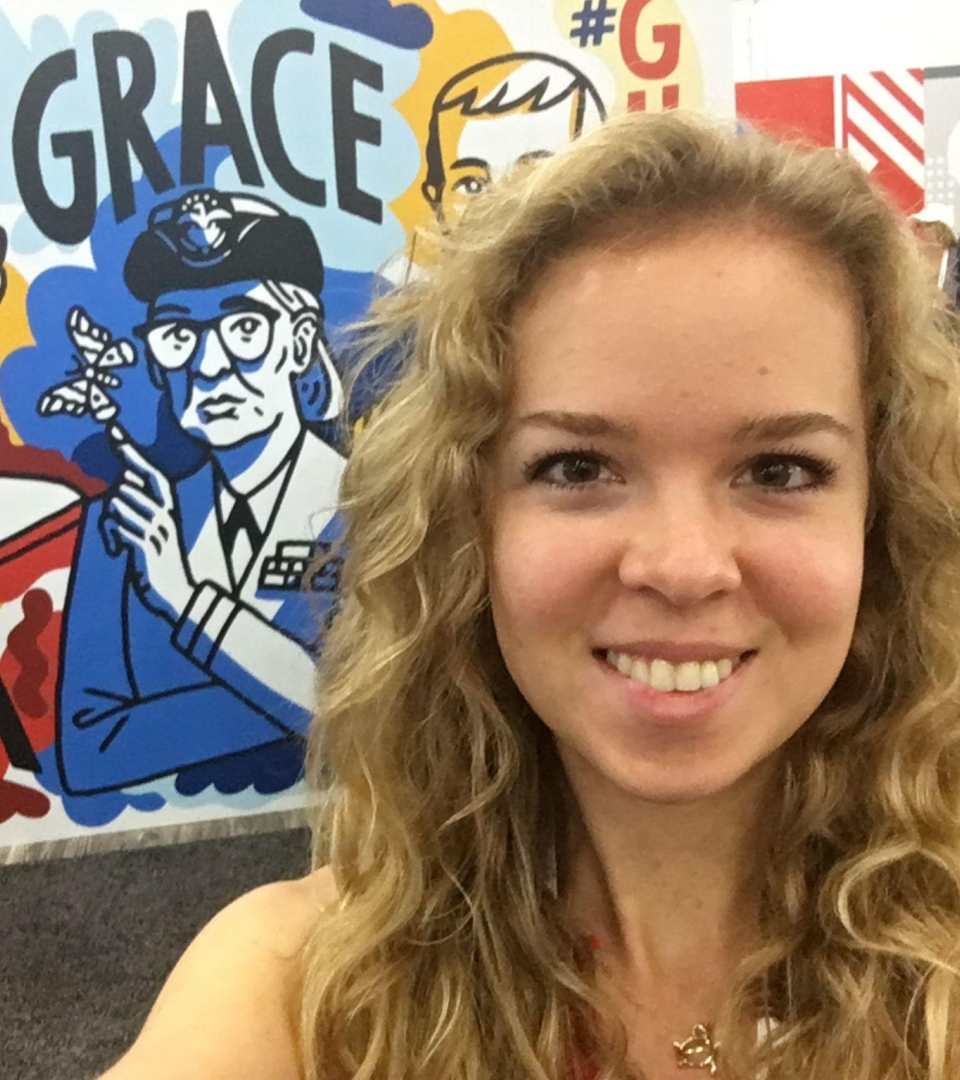# Reverse EXPO 2022

# Hello!

**I'm Andrea**, a Master's student in Computing Science.

My primary research area is NLP, specifically in applying **data science** and **machine learning** to **social communication**.

# Framing Climate Change

"To frame is to select some **aspects** of a perceived reality and make them more **salient** in a communicating text, in such a way as to promote a particular **problem** definition, **causal** interpretation, **moral** evaluation, and/or **treatment** recommendation for the item described."

- frames are manifested through **frame elements**: *keywords*, *stock phrases*, stereotyped images, *sources of information*, and *sentences*

- receivers' responses are clearly affected if they perceive and process information about *one interpretation* and possess little or incommensurable data about *alternatives*

# Framing Climate Change

- How are climate change issues framed on social media?

- **Problem**: no scalable model of general and topical issue frames, nor a flexible methodology for creating one

- **Approach**: using ADaPT-ML to implement frame elements as labelling functions



Electroverse @Electroversenet · 13h

HISTORICAL DATA DESTROYS THE GLOBAL WARMING MYTH, & PEOPLE ARE WAKING TO IT

Of the 50 U.S. state record high temps, 23 were set during the 1930s, while 36 occurred prior to 1960 -- #climatechange proponents are feeding us a fairy tale, and I'm sick of it.

21        379        549



ltfc Australia @ltfc_australia · 13h

So far I've been at fires in nsw and Victoria, in forrests, out on grassland and even in rain forests... if anyone still think this is normal and has nothing to do with #climatechange please unfollow me now. I'm done with you #AustraliaOnFire

39        374        1K

4

# ADaPT-ML

A Data Programming Template for Machine Learning

# Overview

1. Background

2. Why make ADaPT-ML? What is it?

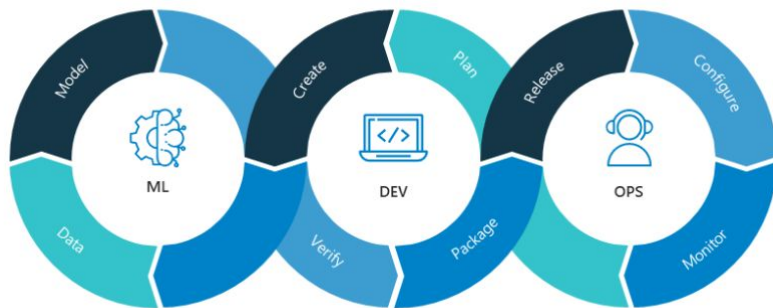3. System walkthrough

4. ADaPTing to new use cases

# Background

**Machine Learning (ML)**

- Machine learning algorithms build a **model** based on sample data, known as **training data**, in order to make **predictions** or decisions without being explicitly programmed to do so.
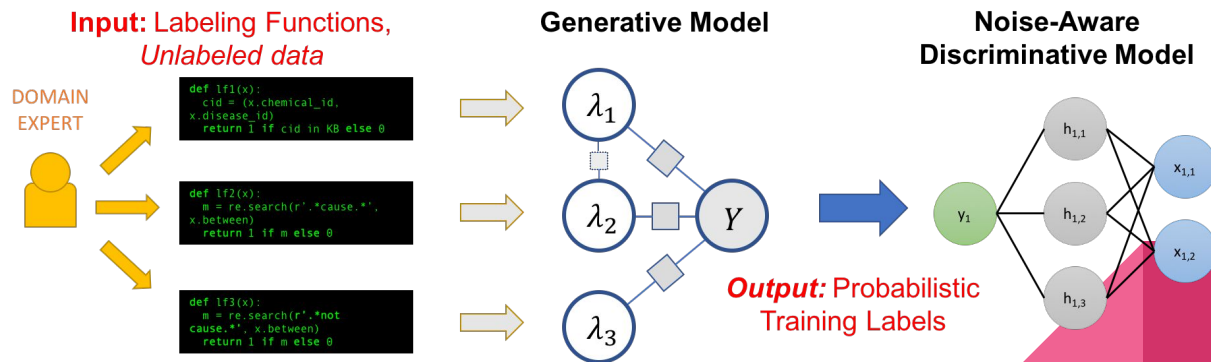
**Machine Learning Operations (MLOps)**

- a set of practices that aims to **deploy** and maintain machine learning **models** in production **reliably and efficiently**.

# Background

## Data Programming

- "a paradigm for the programmatic **creation of training sets** called data programming in which users express **weak supervision strategies** or domain **heuristics** as **labeling functions**, which are programs that label subsets of the data, but that are noisy and may conflict. We show that by explicitly representing this training set labeling process as a **generative model**, we can "**denoise**" the generated training set"

# Why make ADaPT-ML?

Have you needed to model a phenomenon for which there is **insufficient training data**, and **no resources** to acquire it?

> *I have!*

Is there an **open-source** MLOps platform out there with **data programming** at its core?

> *There was not!*

Could you manage by simply using **Snorkel** as a **standalone tool**?

> *You could, but we all know that not having full integration causes headaches.*

# So, what exactly is ADaPT-ML?

It is a **multimodal**-ready **MLOps** system that covers the data processing, **data labelling**, model design, model training and optimization, and endpoint deployment.

This software was created especially for any researcher with:

- Some programming experience or interest in learning how to write code based off of examples.

- Access to large amounts of unlabelled data that is constantly changing, such as social media data.

- Domain expertise or an intuition about how they would follow rules, heuristics, or use knowledge bases to annotate the unlabelled data.

# Questions?

# System Walkthrough

## Data Prerequisites

- Large amounts of unlabelled data in CrateDB

- Features extracted for labelling functions

- Feature vectors for machine learning

# System Walkthrough

## Create a Human-annotated Testing Dataset

- Sample some data points to manually label in Label Studio

- Customize the labelling configuration to handle images, text, audio, conversations, timelines, etc.

- Export the annotations and determine the inter-annotator agreement

- Choose how you want to handle annotator disagreements to finalize one set of gold labels

# System Walkthrough

## Programmatically Create Training Data

- Sample some data points from CrateDB that represent your domain

- Run an experiment with a set of base and custom parameters, including your labelling function features

  - The testing dataset is used to validate the label model
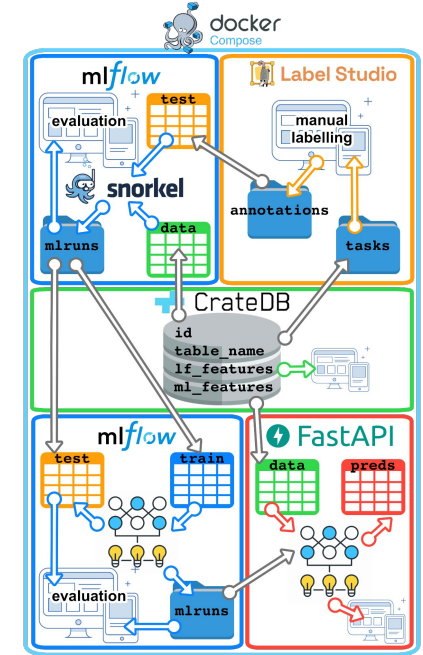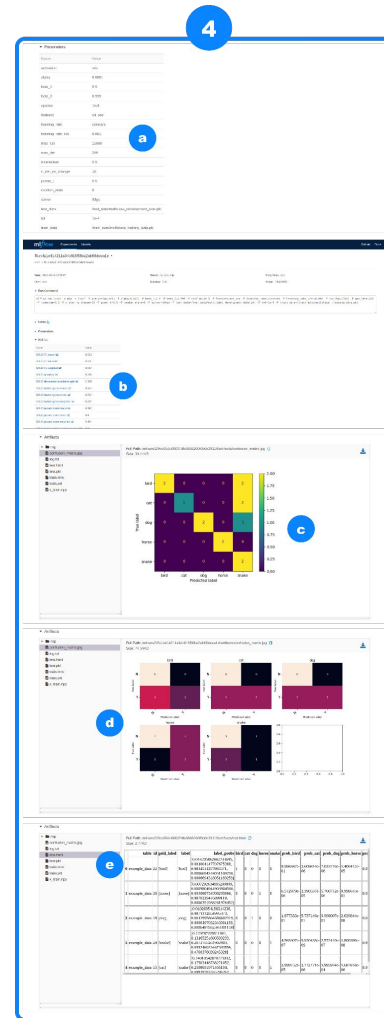
- Monitor the results with the MLflow UI

# System Walkthrough

## Create the Discriminative End Model

- Run an experiment using a suitable machine learning algorithm and your feature vectors

  - The model will be evaluated using the gold labels and label model's labels

- Observe the results with the MLflow UI



amwhitta@ualberta.ca

15

# System Walkthrough

**Deploy the Best Model**

- Get predictions for any data point, needing only to supply the table name and id

- Can keep track of which model has been deployed within the MLflow UI

# Questions?

amwhitta@ualberta.ca

# ADaPTing to New Use Cases          * = future features

**Label Studio**

- Define classification task name and categories

- Format the Labelling Configuration


- ★ Implement methods for taking representative samples

- ★ Expand data formatting for more modalities



Label Every Data Type

Images | Audio | Text | Time Series | Multi-Domain

**Computer Vision**

**Image Classification**
Put images into categories

**Object Detection**
Detect objects on image, bboxes, polygons, circular, and keypoints supported

**Semantic Segmentation**
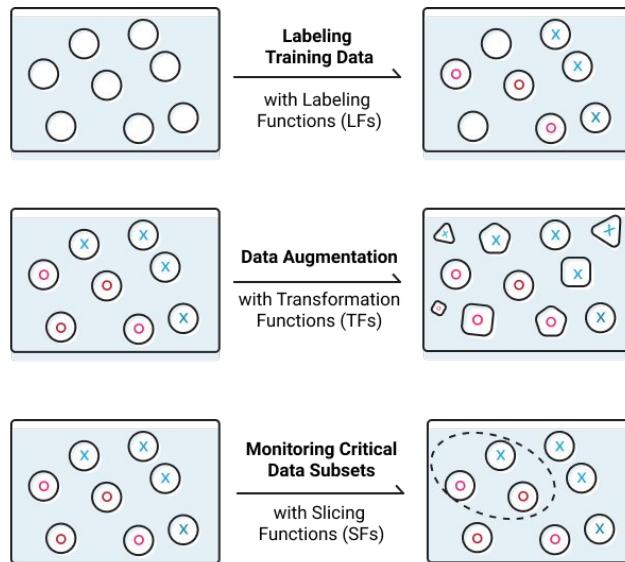Partition image into multiple segments. Use ML models to pre-label and optimize the process

# ADaPTing to New Use Cases       * = future features

**Data Programming**

- Write your labelling functions

- Create an MLflow endpoint

  - parameters for labelling functions

  - loader functions for labelling function features

★ Integrate slicing and transformation functions



Labeling Training Data with Labeling Functions (LFs)

Data Augmentation with Transformation Functions (TFs)

Monitoring Critical Data Subsets with Slicing Functions (SFs)

# ADaPTing to New Use Cases          * = future features

**Model Creation and Deployment**

- Define the response format and endpoint

- ★  Implement more machine learning algorithms

**Thank you for your interest in this project!**

https://github.com/U-Alberta/ADaPT-ML

amwhitta@ualberta.ca

# Resources

https://academic.oup.com/joc/article/43/4/51-58/4160153                 (Entman's publication on framing)
https://blogs.nvidia.com/blog/2020/09/03/what-is-mlops/              (MLOps diagram)
https://en.wikipedia.org/wiki/Machine_learning                     (definition of ML)
https://en.wikipedia.org/wiki/MLOps                            (definition of MLOps)
https://arxiv.org/abs/1605.07723                            (publication describing data programming)
https://www.snorkel.org/blog/snorkel-programming        (diagram and description of Snorkel)
https://labelstud.io/                                       (Label Studio's website)
https://www.snorkel.org/features/                          (more description of Snorkel)